# TCSS 422: OPERATING SYSTEMS

## Multi-level Feedback Queue II, Proportional Share Schedulers, Introduction to Concurrency

Wes J. Lloyd
School of Engineering and Technology
University of Washington - Tacoma

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma

1

---

## TEXT BOOK COUPON

- 15% off textbook code: **LULUBOOKS15**
  (*through Friday Apr 19*)

- https://www.lulu.com/shop/andrea-arpaci-dusseau-and-remzi-arpaci-dusseau/operating-systems-three-easy-pieces-hardcover-version-110/hardcover/product-15gjeeky.html?q=three+easy+pieces+operating+systems&page=1&pageSize=4

- With coupon textbook is only $33.79 + tax & shipping

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.2

2

---

## OFFICE HOURS – SPRING 2024

- **Tuesdays after class until 7:00pm**
  Hybrid (In-person/Zoom)
  - This session will be in person in CP 229.
  - Zoom will be monitored when no student is in CP 229.
- Thursdays after class until 7:00pm – Hybrid (In-person/Zoom)
  - Additional office time will be held on Thursdays after class when there is high demand indicated by a busy Tuesday office hour
  - When Thursday Office Hours are planned, Zoom links will be shared via Canvas
  - Questions after class on Thursdays are always entertained even when the formal office hour is not scheduled

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.3

3

---

## TCSS 422 DISCORD SERVER

- Please join the TCSS 422 A – Spring 2024 Discord Server

- https://discord.gg/H7PPZ5ArFW

- Under Edit Server Profile:
  Please update your 'Server Nickname'
  to your real name or UW NET ID
  THANK YOU

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
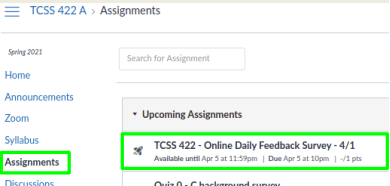L7.4

4

---

## OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.5

5

---

## ONLINE DAILY FEEDBACK SURVEY

- Daily Feedback Quiz in Canvas – Available After Each Class
- Extra credit available for completing surveys *ON TIME*
- Tuesday surveys: due by ~ Wed @ 11:59p
- Thursday surveys: due ~ Mon @ 11:59p

TCSS 422 A › Assignments

Spring 2021

Home
Announcements
Zoom
Syllabus
Assignments
Discussions

Search for Assignment

▾ Upcoming Assignments

TCSS 422 - Online Daily Feedback Survey - 4/1
Available until Apr 5 at 11:59pm  |  Due Apr 5 at 10pm  |  -/1 pts

April 16, 2024
TCSS422: Computer Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.6

6

---

## Slide 7

TCSS 422 - Online Daily Feedback Survey - 4/1

Quiz Instructions

☐ **Question 1**                                                          0.5 pts

On a scale of 1 to 10, please classify your perspective on material covered in today's class:

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |

Mostly
Review To Me          Equal
New and Review                    Mostly
New to Me

☐ **Question 2**                                                          0.5 pts

Please rate the pace of today's class:

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |

Slow          Just Right          Fast

April 16, 2024      TCSS422: Computer Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma      L7.7

---

## Slide 8

# MATERIAL / PACE

- Please classify your perspective on material covered in today's class (31 respondents):
- 1-mostly review, 5-equal new/review, 10-mostly new
- **Average – 6.81  (↓ - previous 7.00)**

- Please rate the pace of today's class:
- 1-slow, 5-just right, 10-fast
- **Average – 5.42 (↑ - previous 5.21)**

April 16, 2024      TCSS422: Computer Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma      L7.8

---

## Slide 9

# FEEDBACK FROM 4/11

- *In the MLFQ, why does job A in Q0 take longer than the 10 ms time slice?*
- This question relates to the teaser problem at the end of Chapter 8.
- Some combination of n short jobs runs for a total of 10 ms per cycle without relinquishing the CPU
  - E.g. 2 jobs = 5 ms ea; 3 jobs = 3.33 ms ea, 10 jobs = 1 ms ea
  - n jobs always uses full time quantum in highest queue (10 ms)
  - Batch job 'A' starts, runs for full quantum of 10ms, pushed to lower queue
  - All other jobs run in top-most queue and context switch and use up 10ms so that job 'A' is starved after it initially runs once
  - If 10ms goal is 5% of the CPU, when must the priority boost be ???
  - **ANSWER → *Priority boost should occur every 200ms***

April 16, 2024      TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma      L7.9

---

## Slide 10

# FEEDBACK - 2

- *How can an MLFQ scheduler account for user advice?*
- We will learn about the Linux '`nice`' command in Ch. 9
- When user runs the nice command it provides a suggestion to the OS scheduler to increase or decrease process priority
- User processes have a maximum priority, and OS processes can still obtain priority greater than any user process priority

- **TRUE/FALSE Question:**
- Round robin is the scheduler that best addresses fairness and average response time of jobs.

April 16, 2024      TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma      L7.10

---

## Slide 11

# OBJECTIVES – 4/16

- Questions from 4/11
- **Assignment 0 - Due Fri Apr 19**
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

April 16, 2024      TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma      L7.11

---

## Slide 12

# ASSIGNMENT 0 - DUE FRI APR 19

- Due Friday April 19 @ 11:59pm
- Grace period: submission ok until Sun Apr 21 @ 11:59 **PM**
- Late submissions thru Tuesday Apr 23 @ 11:59pm

April 16, 2024      TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma      L7.12

## Slide 13

### OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- **C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26**
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.13

13

## Slide 14

### OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- **Quiz 1 and Quiz 2**
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.14

14

## Slide 15

### QUIZ 1

- Active reading on Chapter 9 – Proportional Share Schedulers

- Posted in Canvas
- Due Thursday April 25th at 11:59pm

- **Link:**
- **https://faculty.washington.edu/wlloyd/courses/tcss422/quiz/TCSS422_s2024_quiz_1.pdf**

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.15

15

## Slide 16

### QUIZ 2

- Canvas Quiz – Practice CPU Scheduling Problems

- Posted in Canvas
- Unlimited attempts permitted
- Due Tuesday April 30th at 11:59pm

- **Link:**
- **https://canvas.uw.edu/courses/1728244/quizzes/2030525**

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.16

16

## Slide 17

### COMING SOON...

- Assignment #1
  - To be posted soon

- Midterm Exam
  - Thursday May 2nd
  - In Class

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.17

17

## Slide 18

### OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - **Examples**
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.18

18

## Slide 19

Jackson deploys a 3-level MLFQ scheduler. The time slice is 1 for high priority jobs, 2 for medium priority, and 4 for low priority. This MLFQ scheduler performs a Priority Boost every 6 timer units. When the priority boost fires, the current job is preempted, and the next scheduled job is run in round-robin order.

**SANITY CHECK: Consider the timing graph x-axis should not exceed the combined job length of all jobs.**

| Job | Arrival Time | Job Length |
|-----|--------------|------------|
| A   | T=0          |            |
| B   | T=0          |            |
| C   | T=0          |            |

(11 points) Show a scheduling graph for the MLFQ scheduler for the jobs above. Draw vertical lines for key events and be sure to label the X-axis times as in the example. Please draw clearly. An unreadable graph will loose points.



## Slide 20

Jackson deploys a 3-level MLFQ scheduler. The time slice is 1 for high priority jobs, 2 for medium priority, and 4 for low priority. This MLFQ scheduler performs a Priority Boost every 6 timer units. When the priority boost fires, the current job is preempted, and the next scheduled job is run in round-robin order.

| Job | Arrival Time | Job Length |
|-----|--------------|------------|
| A   | T=0          | 4          |
| B   | T=0          | 16         |
| C   | T=0          | 8          |

time slice is job time

(11 points) Show a scheduling graph for the MLFQ scheduler for the jobs above. Draw vertical lines for key events and be sure to label the X-axis times as in the example. Please draw clearly. An unreadable graph will loose points.



## Slide 21

### EXAMPLE

- Question:
- Given a system with a quantum length of 10 ms *for all jobs* in its highest queue, how often would you have to boost job A (the first job to arrive and run) back to the highest priority level to guarantee that job A, a long-running (and potentially starving) job gets at least 5% of the CPU assuming that on priority boost job execution resets to the front of the queue?

## Slide 22

### EXAMPLE

- Question:
- Given a system with a quantum length of 10 ms *for all jobs* in its highest queue, how often would you have to boost job A (the first job to arrive and run) back to the highest priority level to guarantee that job A, a long-running (and potentially starving) job gets at least 5% of the CPU assuming that on priority boost job execution resets to the front of the queue?

## Slide 23

### EXAMPLE

- Question:
- Given a system with a quantum length of 10 ms *for all jobs* in its highest queue, how often would you have to boost jobs back to the highest priority level to guarantee that a single long-running (and potentially starving) job gets at least 5% of the CPU?

- Some combination of n short jobs runs for a total of 10 ms per cycle without relinquishing the CPU
  - E.g. 2 jobs = 5 ms ea; 3 jobs = 3.33 ms ea, 10 jobs = 1 ms ea
  - n jobs always uses full time quantum in highest queue (10 ms)
  - Batch jobs starts, runs for full quantum of 10ms, pushed to lower queue
  - All other jobs run and context switch totaling the quantum per cycle
  - If 10ms is 5% of the CPU, when must the priority boost be ???
  - **ANSWER → *Priority boost should occur every 200ms***

## Slide 24

### OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - **Lottery scheduler**
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

## Slide 25

# CHAPTER 9 - PROPORTIONAL SHARE SCHEDULER

April 16, 2024 — TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma

L7.25

25

## Slide 26

# PROPORTIONAL SHARE SCHEDULER

- Also called fair-share scheduler or lottery scheduler
  - Guarantees each job receives some percentage of CPU time based on share of "tickets"
  - Each job receives an allotment of tickets
  - % of tickets corresponds to potential share of a resource
  - Can conceptually schedule any resource this way
    - CPU, disk I/O, memory

April 16, 2024 — TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma

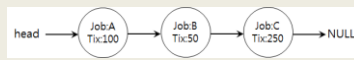L7.26

26

## Slide 27

# LOTTERY SCHEDULER

- Simple implementation
  - Just need a random number generator
    - Picks the winning ticket
  - Maintain a data structure of jobs and tickets (list)
  - Traverse list to find the owner of the ticket
  - Consider sorting the list for speed

April 16, 2024 — TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma

L7.27

27

## Slide 28

# LOTTERY SCHEDULER IMPLEMENTATION



```
1    // counter: used to track if we've found the winner yet
2    int counter = 0;
3
4    // winner: use some call to a random number generator to
5    // get a value, between 0 and the total # of tickets
6    int winner = getrandom(0, totaltickets);
7
8    // current: use this to walk through the list of jobs
9    node_t *current = head;
10
11   // loop until the sum of ticket values is > the winner
12   while (current) {
13           counter = counter + current->tickets;
14           if (counter > winner)
15                   break; // found the winner
16           current = current->next;
17   }
18   // 'current' is the winner: schedule it...
```

April 16, 2024 — TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma

L7.28

28

## Slide 29

# OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

April 16, 2024 — TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma

L7.29

29

## Slide 30

# TICKET MECHANISMS

- Ticket currency / exchange
  - User allocates tickets in any desired way
  - OS converts user currency into global currency

- Example:
  - There are 200 global tickets assigned by the OS

  User A → 500 (A's currency) to A1 → 50 (global currency)
          → 500 (A's currency) to A2 → 50 (global currency)

  User B → 10 (B's currency) to B1 → 100 (global currency)

April 16, 2024 — TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma

L7.30

30

## TICKET MECHANISMS - 2

- Ticket transfer
  - Temporarily hand off tickets to another process

- Ticket inflation
  - Process can temporarily raise or lower the number of tickets it owns
  - If a process needs more CPU time, it can boost tickets.

31

## LOTTERY SCHEDULING

- Scheduler picks a **winning** ticket
  - Load the job with the winning ticket and run it

- Example:
  - Given 100 tickets in the pool
    - Job A has 75 tickets: 0 - 74
    - Job B has 25 tickets: 75 – 99

| Scheduler's winning tickets: | 63 85 70 39 76 17 29 41 36 39 10 99 68 83 63 |
| --- | --- |
| Scheduled job: | A  B  A  A  B  A  A  A  A  A  B  A  B  A |

- But what do we know about probability of a coin flip?

32

## COIN FLIPPING

- Equality of distribution (fairness) requires a lot of flips!



Similarly,
Lottery scheduling requires lots of "rounds" to achieve fairness.

33

## LOTTERY FAIRNESS

- With two jobs
  - Each with the same number of tickets (t=100)



When the job length is not very long,
average unfairness can be **quite severe**.

34

## LOTTERY SCHEDULING CHALLENGES

- What is the best approach to assign tickets to jobs?
  - Typical approach is to assume users know best
  - Users are provided with tickets, which they allocate as desired

- How should the OS automatically distribute tickets upon job arrival?
  - What do we know about incoming jobs a priori ?
  - Ticket assignment is really an open problem…

35

# WE WILL RETURN AT 4:50PM

36

## OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - **Stride scheduler**
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.37

37

## STRIDE SCHEDULER

- Addresses statistical probability issues with lottery scheduling

- Instead of guessing a random number to select a job, simply count…

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.38

38

## STRIDE SCHEDULER - 2

- Jobs have a "stride" value
  - A stride value describes the counter pace when the job should give up the CPU
  - Stride value is **Inverse in proportion** to the job's number of tickets (more tickets = smaller stride)

- Total system tickets = 10,000
  - Job A has 100 tickets → $A_{stride}$ = 10000/100 = 100 stride
  - Job B has 50 tickets → $B_{stride}$ = 10000/50 = 200 stride
  - Job C has 250 tickets → $C_{stride}$ = 10000/250 = 40 stride

- Stride scheduler tracks "pass" values for each job (A, B, C)

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.39

39

## STRIDE SCHEDULER - 3

- Basic algorithm:
  1. Stride scheduler picks job with the lowest pass value
  2. Scheduler increments job's pass value by its stride and starts running
  3. Stride scheduler increments a counter
  4. When counter exceeds pass value of current job, pick a new job (go to 1)

- <u>KEY:</u> When the counter reaches a job's "PASS" value, the scheduler <u>passes</u> on to the next job…

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.40

40

## STRIDE SCHEDULER - EXAMPLE

- Stride values
  - Tickets = priority to select job
  - Stride is inverse to tickets
  - Lower stride = more chances to run <u>(higher priority)</u>

    <u>Priority</u>
    C stride = 40
    A stride = 100
    B stride = 200

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.41

41

## STRIDE SCHEDULER EXAMPLE - 2

- <u>Three-way tie:</u> randomly pick job A (all pass values=0)
- Set A's pass value to A's stride = 100
- Increment counter until > 100
- Pick a new job: <u>two-way tie</u>

**Tickets**
C = 250
A = 100
B = 50

| Pass(A) (stride=100) | Pass(B) (stride=200) | Pass(C) (stride=40) | Who Runs? |
|---|---|---|---|
| 0 | 0 | 0 | A |
| 100 | 0 | 0 | B |
| 100 | 200 | 0 | C |
| 100 | 200 | 40 | C |
| 100 | 200 | 80 | C |
| 100 | 200 | 120 | A |
| 200 | 200 | 120 | C |
| 200 | 200 | 160 | C |
| 200 | 200 | 200 | … |

◄ Initial job selection is random. All @ 0

◄ C has the most tickets and receives a lot of opportunities to run…

April 16, 2024 | TCSS422: Operating Systems [Spring 2024] School of Engineering and Technology, University of Washington - Tacoma | L7.42

42

## Slide 43

### STRIDE SCHEDULER EXAMPLE - 3

- We set A's counter (pass value) to A's stride = 100
- Next scheduling decision between B (pass=0) and C (pass=0)
  - Randomly choose B
- C has the lowest counter for next 3 rounds

**Tickets**
C = 250
A = 100
B = 50

| Pass(A) (stride=100) | Pass(B) (stride=200) | Pass(C) (stride=40) | Who Runs? |
|---|---|---|---|
| 0 | 0 | 0 | A |
| 100 | 0 | 0 | B |
| 100 | 200 | 0 | C |
| 100 | 200 | 40 | C |
| 100 | 200 | 80 | C |
| 100 | 200 | 120 | A |
| 200 | 200 | 120 | C |
| 200 | 200 | 160 | C |
| 200 | 200 | 200 | ... |

← C has the most tickets and is selected to run more often ...

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.43

43

## Slide 44

### STRIDE SCHEDULER EXAMPLE - 4

- Job counters support determining which job to run next
- Over time jobs are scheduled to run based on their priority represented as their **share of tickets...**
- **Tickets are analogous to job priority**

**Tickets**
C = 250
A = 100
B = 50

| Pass(A) (stride=100) | Pass(B) (stride=200) | Pass(C) (stride=40) | Who Runs? |
|---|---|---|---|
| 0 | 0 | 0 | A |
| 100 | 0 | 0 | B |
| 100 | 200 | 0 | C |
| 100 | 200 | 40 | C |
| 100 | 200 | 80 | C |
| 100 | 200 | 120 | A |
| 200 | 200 | 120 | C |
| 200 | 200 | 160 | C |
| 200 | 200 | 200 | ... |

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.44

44

## Slide 45

### Which of the following is NOT a problem with proportional share schedulers?

How tickets should be distributed to incoming jobs — A

Lottery scheduler is only eventually fair — B

Given 2 users A and B who both receive a 50% timeshare of the system, the runtime for User A's jobs is dependent on the runtime of User B's. — C

All of the above — D

None of the above — E

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.45

45

## Slide 46

### OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - **Linux Completely Fair Scheduler**
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.46

46

## Slide 47

### LINUX: COMPLETELY FAIR SCHEDULER (CFS)

- Large Google datacenter study:
  *"Profiling a Warehouse-scale Computer"* (Kanev et al.)
- Monitored 20,000 servers over 3 years
- Found 20% of CPU time spent in the Linux kernel
- 5% of CPU time spent in the CPU scheduler!
- Study highlights importance for high performance OS kernels and CPU schedulers !

Figure 5: Kernel time, especially time spent in the scheduler, is a significant fraction of WSC cycles.

See: https://dl.acm.org/doi/pdf/10.1145/2749469.2750392

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.47

47

## Slide 48

### LINUX: COMPLETELY FAIR SCHEDULER (CFS)

- Loosely based on the stride scheduler
- CFS models system as a Perfect Multi-Tasking System
  - In a perfect system every process of the same priority (class) receives exactly $1/n^{th}$ of the CPU time
- Each scheduling class has a runqueue
  - Groups processes of the same class
  - In the class, scheduler picks task w/ lowest `vruntime` to run
  - Time slice varies based on how many jobs in shared runqueue
  - Minimum time slice prevents too many context switches (e.g. 3 ms)

April 16, 2024
TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma
L7.48

48

## COMPLETELY FAIR SCHEDULER - 2

- Every thread/process has a scheduling class (policy):
- **Normal classes**: SCHED_OTHER (TS), SCHED_IDLE, SCHED_BATCH
  - TS = Time Sharing
- **Real-time classes**: SCHED_FIFO (FF), SCHED_RR (RR)
- How to show scheduling class and priority:
- ```
  #class
  ps –elfc
  ```
- ```
  #priority (nice value)
  ps ax -o pid,ni,cls,pri,cmd
  ```

| April 16, 2024 | TCSS422: Operating Systems [Spring 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.49 |

49

## COMPLETELY FAIR SCHEDULER - 3

- Linux ≥ 2.6.23: Completely Fair Scheduler (CFS)
- Linux < 2.6.23: O(1) scheduler
- Linux maintains simple counter (`vruntime`) to track how long each thread/process has run
- CFS picks process with lowest `vruntime` to run next
- CFS adjusts timeslice based on # of proc waiting for the CPU
- Kernel parameters that specify CFS behavior:
  ```
  $ sudo sysctl kernel.sched_latency_ns
  kernel.sched_latency_ns = 24000000
  $ sudo sysctl kernel.sched_min_granularity_ns
  kernel.sched_min_granularity_ns = 3000000
  $ sudo sysctl kernel.sched_wakeup_granularity_ns
  kernel.sched_wakeup_granularity_ns = 4000000
  ```

| April 16, 2024 | TCSS422: Operating Systems [Spring 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.50 |

50

## COMPLETELY FAIR SCHEDULER - 4

- `Sched_min_granularity_ns` (3ms)
  - Time slice for a process: busy system (w/ full runqueue)
  - If system has idle capacity, time slice exceeds the min as long as difference in `vruntime` between running process and process with lowest `vruntime` is less than `sched_wakeup_granularity_ns` (4ms)
- Scheduling time period is: total cycle time for iterating through a set of processes where each is allowed to run (like round robin)
- Example:
  `sched_latency_ns` (24ms)
  if (proc in runqueue < `sched_latency_ns`/`sched_min_granularity`)
  or
  `sched_min_granularity` * number of processes in runqueue
  Ref: https://www.systutorials.com/sched_min_granularity_ns-sched_latency_ns-cfs-affect-timeslice-processes/

| April 16, 2024 | TCSS422: Operating Systems [Spring 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.51 |

51

## CFS TRADEOFF

- **HIGH**        `sched_min_granularity_ns` (timeslice)
              `sched_latency_ns`
              `sched_wakeup_granularity_ns`

  CFS features reduced context switching → less overhead poor near-term fairness

- **LOW**        `sched_min_granularity_ns` (timeslice)
              `sched_latency_ns`
              `sched_wakreup_granularity_ns`

  CFS features increased context switching → more overhead better near-term fairness

| April 16, 2024 | TCSS422: Operating Systems [Spring 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.52 |

52

## COMPLETELY FAIR SCHEDULER - 5

- Runqueues are stored using a Linux red-black tree
  - Self balancing binary tree - nodes indexed by `vruntime`
- Leftmost node has lowest `vruntime` (approx execution time)
- Walking tree to find left most node has very low big O complexity:
  - *~O(log N) for N nodes*
- Completed processes are removed



Nodes represent sched_entity(s) indexed by their virtual runtime

Virtual runtime

Most need of CPU                                  Least need of CPU

| April 16, 2024 | TCSS422: Operating Systems [Spring 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.53 |

53

## CFS: JOB PRIORITY



- Time slice: Linux *"Nice value"*
  - Nice predates the CFS scheduler
  - Top shows nice values
  - Process command (nice & priority):
    `ps ax -o pid,ni,cmd,%cpu, pri`
- Nice Values: from -20 to 19
  - Lower is **_higher_** priority, default is 0
  - `vruntime` is a weighted time measurement
  - Priority weights the calculation of `vruntime` within a runqueue to give high priority jobs a *boost*.
    - Influences job's position in rb-tree

| April 16, 2024 | TCSS422: Operating Systems [Spring 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.54 |

54

## COMPLETELY FAIR SCHEDULER - 6

- CFS tracks cumulative job run time with the `vruntime` variable
- The task on a given runqueue with the lowest `vruntime` is scheduled next
- `struct sched_entity` contains `vruntime` parameter
  - Describes process execution time in nanoseconds
  - Value is not pure runtime, is weighted based on job priority
  - **GOAL:** Perfect scheduler →
    achieve equal `vruntime` for all processes of same priority
- Sleeping jobs: upon return a temporary `vruntime` can be used to increase temporarily the priority of the task
- When tasks wait for I/O they should receive a comparable share of the CPU as if they were performing compute ops when run again
- Key takeaway:
  *Identifying the next job to schedule is really fast!*

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma | L7.55

55

## COMPLETELY FAIR SCHEDULER - 7

- More information:
- Man page: "man sched" : Describes Linux scheduling API
- http://manpages.ubuntu.com/manpages/bionic/man7/sched.7.html
- https://www.kernel.org/doc/Documentation/scheduler/sched-design-CFS.txt
- https://en.wikipedia.org/wiki/Completely_Fair_Scheduler
- See paper: The Linux Scheduler – a Decade of Wasted Cores
- http://www.ece.ubc.ca/~sasha/papers/eurosys16-final29.pdf

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington – Tacoma | L7.56

56

## OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - Critical section

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma | L7.57
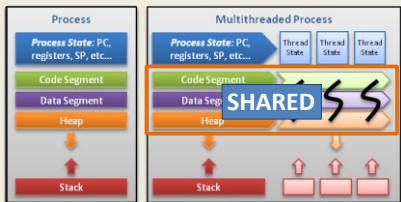
57

## CHAPTER 26 - CONCURRENCY: AN INTRODUCTION

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma | L7.58

58

## THREADS



©Alfred Park, http://randu.org/tutorials/threads

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma | L7.59

59

## THREADS - 2

- Enables a single process (program) to have multiple "workers"
  - This is parallel programming…
- Supports independent path(s) of execution within a program *with shared memory* …
- Each thread has its own Thread Control Block (TCB)
  - PC, registers, SP, and stack
- Threads share code segment, memory, and heap are shared
- *What is an embarrassingly parallel program?*

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma | L7.60

60

## PROCESS AND THREAD METADATA

- Thread Control Block vs. Process Control Block

```
Thread identification              Process identification
Thread state                       Process status
CPU information:                   Process state:
        Program counter                   Process status word
        Register contents                 Register contents
Thread priority                           Main memory
Pointer to process that created this thread   Resources
Pointers to all other threads created by this thread   Process priority
                                   Accounting
```

April 16, 2024 · TCSS422: Operating Systems [Spring 2024] · School of Engineering and Technology, University of Washington - Tacoma · L7.61

61

## SHARED ADDRESS SPACE

- Every thread has it's own stack / PC



A Single-Threaded Address Space / Two threaded Address Space

April 16, 2024 · TCSS422: Operating Systems [Spring 2024] · School of Engineering and Technology, University of Washington - Tacoma · L7.62

62

## THREAD CREATION EXAMPLE

```c
#include <stdio.h>
#include <assert.h>
#include <pthread.h>

void *mythread(void *arg) {
    printf("%s\n", (char *) arg);
    return NULL;
}

int
main(int argc, char *argv[]) {
    pthread_t p1, p2;
    int rc;
    printf("main: begin\n");
    rc = pthread_create(&p1, NULL, mythread, "A"); assert(rc == 0);
    rc = pthread_create(&p2, NULL, mythread, "B"); assert(rc == 0);
    // join waits for the threads to finish
    rc = pthread_join(p1, NULL); assert(rc == 0);
    rc = pthread_join(p2, NULL); assert(rc == 0);
    printf("main: end\n");
    return 0;
}
```

April 16, 2024 · TCSS422: Operating Systems [Spring 2024] · School of Engineering and Technology, University of Washington - Tacoma · L7.63

63

## POSSIBLE ORDERINGS OF EVENTS

| int main() | Thread 1 | Thread 2 |
|---|---|---|
| Starts running | | |
| Prints 'main: begin' | | |
| Creates Thread 1 | | |
| Creates Thread 2 | | |
| Waits for T1 | | |
| | Runs | |
| | Prints 'A' | |
| | Returns | |
| Waits for T2 | | |
| | | Runs |
| | | Prints 'B' |
| | | Returns |
| Prints 'main: end' | | |

April 16, 2024 · TCSS422: Operating Systems [Spring 2024] · School of Engineering and Technology, University of Washington - Tacoma · L7.64

64

## POSSIBLE ORDERINGS OF EVENTS - 2

| int main() | Thread 1 | Thread 2 |
|---|---|---|
| Starts running | | |
| Prints 'main: begin' | | |
| Creates Thread 1 | | |
| | Runs | |
| | Prints 'A' | |
| | Returns | |
| Creates Thread 2 | | |
| | | Runs |
| | | Prints 'B' |
| | | Returns |
| Waits for T1 | Returns immediately | |
| Waits for T2 | | Returns immediately |
| Prints 'main: end' | | |

April 16, 2024 · TCSS422: Operating Systems [Spring 2024] · School of Engineering and Technology, University of Washington - Tacoma · L7.65

65

## POSSIBLE ORDERINGS OF EVENTS - 3

| int main() | Thread 1 | Thread 2 |
|---|---|---|
| Starts running | | |
| Prints 'main: begin' | | |
| Creates Thread 1 | | |
| Creates Thread 2 | | |
| Waits for T1 | | |
| | Runs | |
| | Prints 'A' | |
| | Returns | |
| Waits for T2 | | Immediately returns |
| Prints 'main: end' | | |

**What if execution order of events in the program matters?**

April 16, 2024 · TCSS422: Operating Systems [Spring 2024] · School of Engineering and Technology, University of Washington - Tacoma · L7.66

66

## COUNTER EXAMPLE

- Counter example

- A + B : ordering
- Counter: incrementing global variable by two threads

- *Is the counter example embarrassingly parallel?*

- *What does the parallel counter program require?*

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma | L7.67

67

## PROCESSES VS. THREADS

- What's the difference between forks and threads?
  - Forks: duplicate a process
  - Think of *CLONING* - There will be two identical processes at the end
  - Threads: no duplication of code/heap, lightweight execution threads



April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma | L7.68

68

## OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - **Race condition**
  - Critical section

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington  - Tacoma | L7.69

69

## RACE CONDITION

- What is happening with our counter?
  - When counter=50, consider code: counter = counter + 1
  - If synchronized, counter will = 52

```
                                                  (after instruction)
OS          Thread1         Thread2           PC  %eax  counter
            before critical section           100   0    50
            mov 0x8049a1c, %eax               105  50    50
            add $0x1, %eax                     108  51    50
interrupt
save T1's state
restore T2's state                            100   0    50
                            mov 0x8049a1c, %eax 105  50    50
                            add $0x1, %eax      108  51    50
                            mov %eax, 0x8049a1c 113  51    51
interrupt
save T2's state
restore T1's state                            108  51    50
            mov %eax, 0x8049a1c               113  51    51
```

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma | L7.70

70

## OBJECTIVES – 4/16

- Questions from 4/11
- Assignment 0 - Due Fri Apr 19
- C Tutorial - Pointers, Strings, Exec in C - Due Fri Apr 26
- Quiz 1 and Quiz 2
- Chapter 8: Multi-level Feedback Queue
  - Gaming the Scheduler
  - Examples
- Chapter 9: Proportional Share Schedulers
  - Lottery scheduler
  - Ticket mechanisms
  - Stride scheduler
  - Linux Completely Fair Scheduler
- Chapter 26: Concurrency: An Introduction
  - Introduction
  - Race condition
  - **Critical section**

April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington  - Tacoma | L7.71

71

## CRITICAL SECTION

- Code that accesses a shared variable must not be *concurrently* executed by more than one thread

- Multiple *active* threads inside a *critical section* produce a *race condition*.

- *Atomic execution* (*all code executed as a unit*) must be ensured in *critical* sections
  - These sections must be *mutually exclusive*



April 16, 2024 | TCSS422: Operating Systems [Spring 2024]
School of Engineering and Technology, University of Washington - Tacoma | L7.72

72

## LOCKS

- To demonstrate how critical section(s) can be executed "atomically-*as a unit*" Chapter 27 & beyond introduce locks

```
1   lock_t mutex;
2   . . .
3   lock(&mutex);
4   balance = balance + 1;       Critical section
5   unlock(&mutex);
```

- Counter example revisited

| April 16, 2024 | TCSS422: Operating Systems [Spring 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L7.73 |

73

## QUESTIONS

74